

## STATISTICS

→ Statistics → Statistics is the science of collecting, organizing and analyzing the data.

Data → facts or pieces of information.

## STATISTICS

↓ | ↓  
Descriptive Stats (EDA + FE)

- ① It consists of organizing and summarizing the data.

FE → Feature Engineering

Inferential Stats

It consists of collecting sample data and making conclusion about population data using some experiments.

## Sampling Techniques

- 1) Simple Random Sampling - Every member of the population ( $N$ ) has an equal chance of being selected for your sample.  
Eg:- Exit Poll, General survey, Movie reviews

- 2) Stratified Sampling:-  
Eg:- Based on Gender, Blood Group.  
Strata(layers) → Clusters → Groups

- 3) Systematic Sampling:- Select every  $n^{th}$  individual out of population ( $N$ ).

- 4) Convenient Sampling - Only those who are interested in the survey will only participate.

- \* Variable → A variable is a property that can take any values

i) Quantitative Variable → Measured Numerically [mathematical operations can be done]. Eg- Age, weight, height, rainfall (cm), temp, distance.

ii) Qualitative variable = categorical variables {Based on some characteristics they are grouped together}

Eg:- Gender, Types of flowers, Types of Movies.

### Quantitative Variable

↓  
Discrete Variable  
Eg:- Whole numbers only

Nb. of bank accounts, children in family  
Pincode is an example. Also there can be many number of data points.

↓  
Continuous Variable  
Decimal numbers

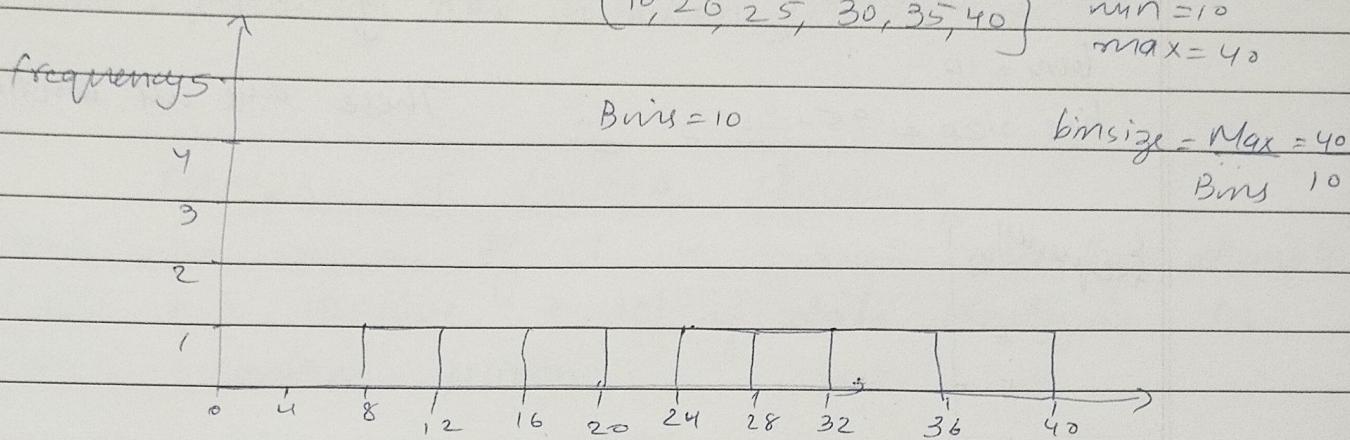
Eg:- Height, weight, ages  
rainfall, speed

### Histograms

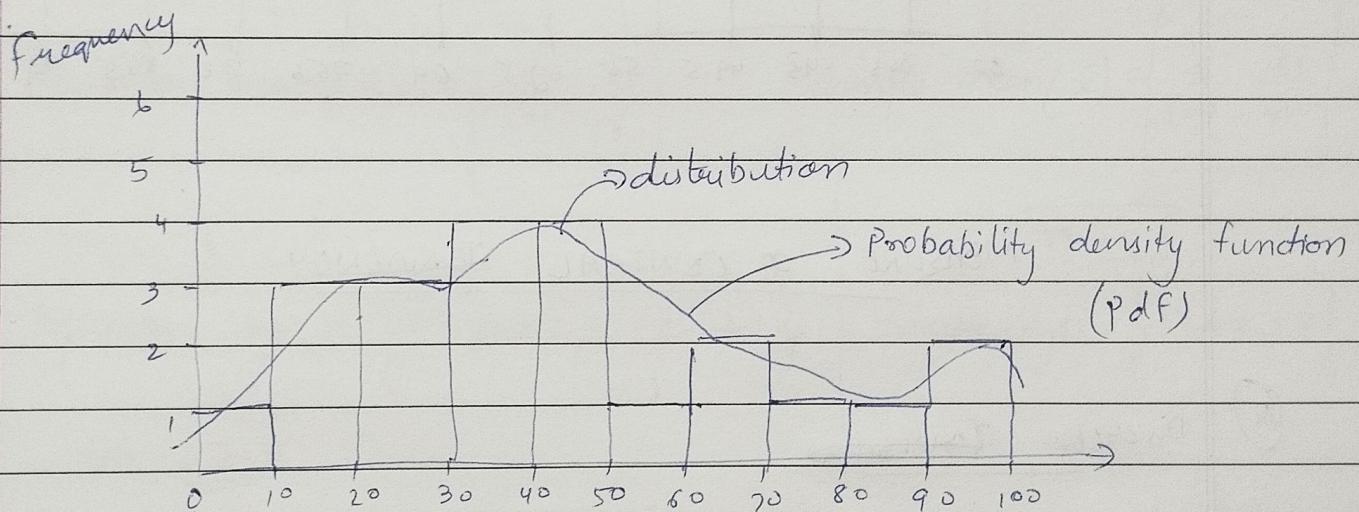
Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

- 1) Sort the numbers
- 2) Bin → No. of groups
- 3) Binsize.

lets, look at an example.

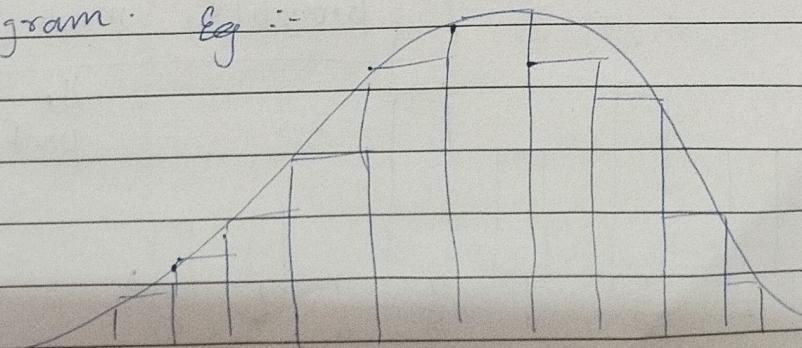


\* for Ages.  $\text{Binsize} = \frac{100}{10} = 10$   
10 → Assumed



\* We are smoothening the curve, this shows how the data is getting distributed.

\* When we smoothen the curve internally there will be histogram. Eg :-

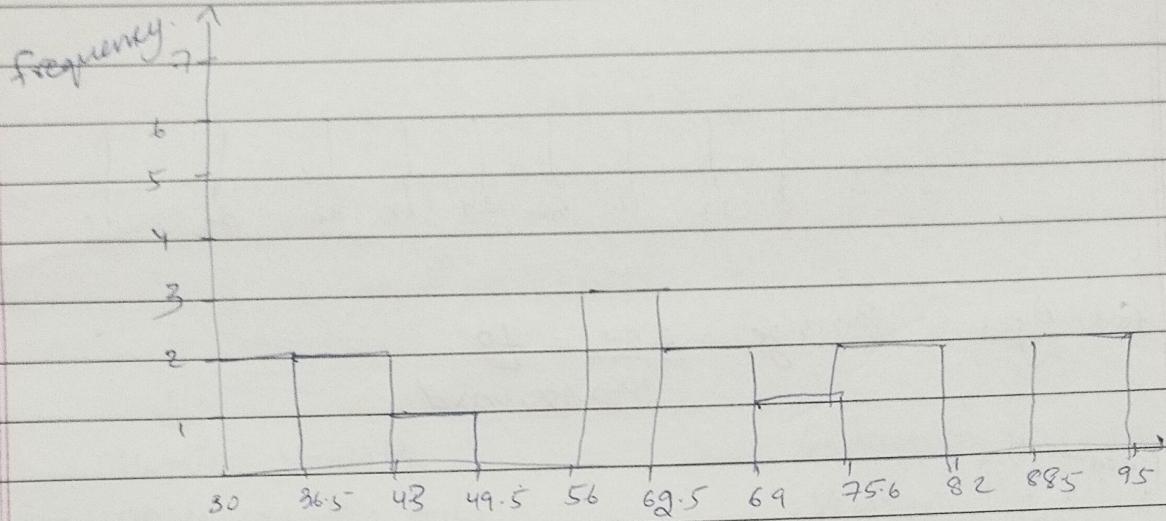


Weight = {30, 35, 38, 42, 46, 58, 67, 62, 63, 68, 75, 77, 80, 90, 95}

lendi = 10

$$\text{bin size} = \frac{95 - 30}{10} = 6.5$$

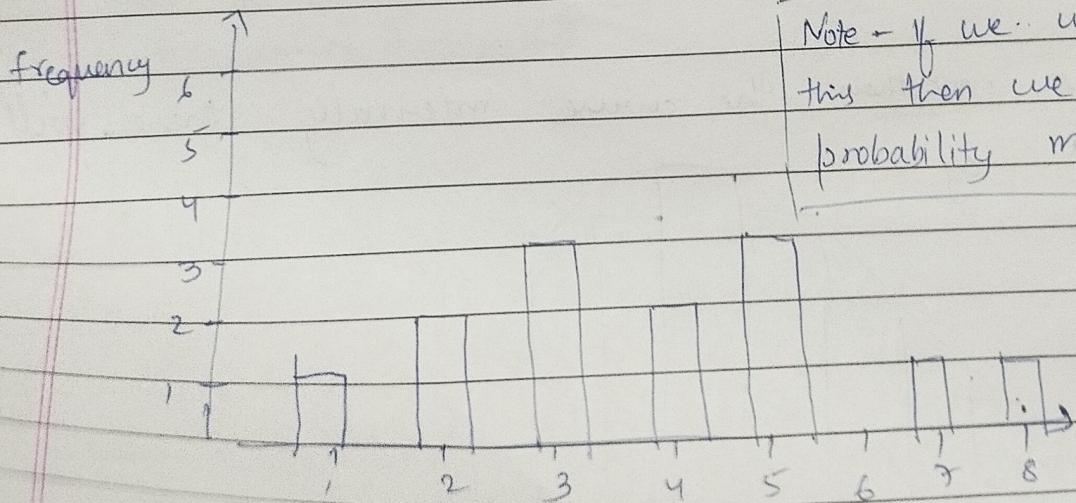
These are all continuous values



### MEASURE OF CENTRAL TENDENCY

Q) Discrete ~~continuous~~

No. of bank accounts = [2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5]



Note - If we will smoothen this then we have to use probability mass function

pmf

pdf : probability density function  $\rightarrow$  continuous  
 pmf : probability mass function  $\int \rightarrow$  discrete

### MEASURE OF CENTRAL TENDENCY

- (1) Mean { A measure of CT is a single value that attempts to describe a set of data identifying the central position.
- (2) Median { to describe a set of data identifying the central position.
- (3) Mode { to describe a set of data identifying the central position.

$$\text{Mean } \bar{x} = \{1, 2, 3, 4, 5\}, \text{ Average/Mean} = \frac{1+2+3+4+5}{5} = 3$$

Population ( $N$ )

$N \geq n$

Sample ( $n$ )

$$\text{Population mean}(\mu) = \left[ \sum_{i=1}^N x_i \right] / N$$

$$\text{Sample mean}(\bar{x}) = \left[ \sum_{i=1}^n x_i \right] / n$$

No 6

$\bar{x} = 4$

$$\text{Population Age} = \{24, 23, 21, 28, 27\} \quad \text{Sample Ages} = \{24, 21, 27\}$$

$$\underline{\underline{\mu = 24.5}}$$

$$\underline{\underline{\bar{x} = 21.5}}$$

$$\begin{array}{|c|c|} \hline \mu & \geq \bar{x} \\ \hline \bar{x} & \geq \mu \\ \hline \end{array}$$

### Practical Application (Feature Engineering)

Age	Salary	family size
-	-	-

nan

-

-

We can use mean in nan.

nan

-

② Median

$$\{1, 2, 3, 4, 5\} \rightarrow \{1, 2, 3, 4, 5, \boxed{100}\}$$

↑ Outlier

$$\bar{x} = 3 \quad \longrightarrow \quad \bar{x} = 19.16$$

Steps to find out median

① Sort the numbers

② find the central number

i) if the no. of elements are even we find the average of central elements.

ii) if the no. of elements are odd we find the central elements.

Sorted

$$\{1, 2, 3, 4, 5, 6, 7, 8, 100, 120\}$$

$$\text{Median} = \underline{\underline{5.5}} = \frac{5+6}{2}$$

$$\text{Mean} = \underline{\underline{25.6}}$$

③ Mode = Most frequent occurring element.

$$\{1, 2, 2, 3, 3, 3, 4, 5\} \rightarrow \text{Mode} = 3$$

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5\} \rightarrow \text{Mode} = [2, 3]$$

\* Measure of Dispersion

① Variance ( $\sigma^2$ ) ← Spread of Data

② Standard deviation ( $\sigma$ )

Variance

Population variance ( $\sigma^2$ )

$$\sigma^2 = \frac{N}{\sum_{i=1}^N} \frac{(x_i - \mu)^2}{N}$$

$$\{1, 2, 3, 4, 5\}$$

$$\mu = 3$$

Sample variance ( $s^2$ )

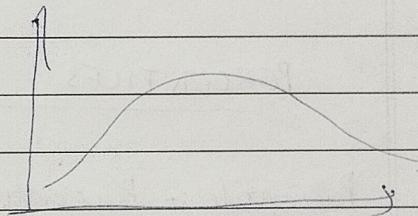
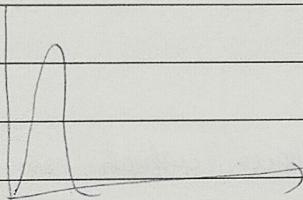
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

$$\{1, 2, 3, 4, 5, 6, 80\}$$

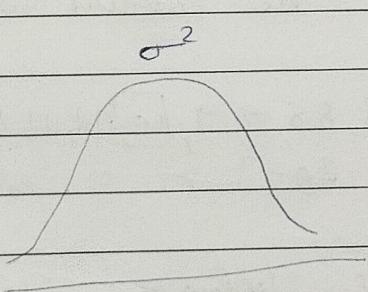
$$\bar{x} = 14.4$$

$$\sigma^2 = 2$$

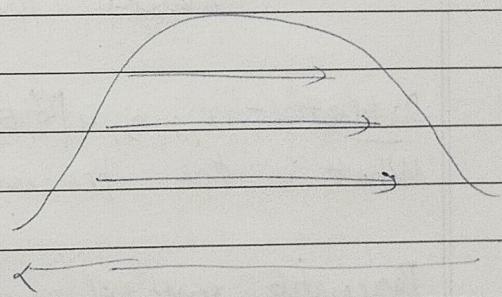
$$\sigma^2 = 11$$



Note - As  $\sigma^2 \uparrow$  spread also  $\uparrow$



L



Q. Why sample variance is divided by  $n-1$ .

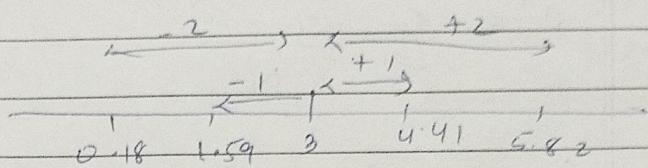
## \* Standard deviation ( $\sqrt{\sigma^2}$ )

$$\{1, 2, 3, 4, 5\}$$

$$\bar{x} = 3$$

$$\sigma^2 = 2$$

$$\sigma = 1.41$$



Standard deviation indicates how many standard deviation the data points are away from the mean.

## PERCENTILES AND QUARTILES

Percentile :- Percentile is a value below which a certain percentage of observation lie.

99 percentile  $\Rightarrow$  It means the person has got better marks than 99% of the entire students.

Dataset :- 2, 2, 3, 4, ~~5~~, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

What is the percentile rank of 10.

$$\text{Percentile rank of } x = \frac{\# \text{ No. of value below } x}{n} = \frac{16}{20} = \frac{4}{5} = 80\%$$

\* What is the value that exists at 25 percentile.

$$0.25 = \frac{n}{20}$$

$n = 5^{\text{th}}$  index (5 is the value)

$\underline{m}$  in dataset consider lower o.

$$\boxed{\text{Value} = \frac{\text{Percentile}}{100} \times n}$$

\* 5 Number summary.

- 1) Minimum
- 2) First Quartile (25 percentile)  $Q_1$
- 3) Median
- 4) Third Quartile (75 percentile)  $Q_3$
- 5) Maximum.

$\{ 1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27 \}$

[Lower fence  $\longleftrightarrow$  Higher fence]

75% 25%

$$IQR = Q_3 - Q_1$$

$$\text{Lower Fence} = Q_1 - 1.5 \text{ (IQR)}$$

$$\text{Higher Fence} = Q_3 + 1.5 \text{ (IQR)}$$

Inter Quartile Range

1.5 Standard Deviation

$$Q_1 = \frac{25 \times (n+1)}{100} = 5.25 \text{ index}$$

Calculate Avg of 5th & 6th index is  $\boxed{3}$  because  $5=25$   
index is not there.

$$Q_3 = \frac{75}{100} (n+1) = 0.75 \times 21 = 15.75 \text{ index}$$

Calculating avg of 15th & 16th index  $\rightarrow \boxed{17.5}$

$$Q_1 = 3$$

$$IQR = 4.5$$

$$Q_3 = 17.5$$

$$\text{Lower fence} = 3 - (1.5) \times 4.5 = -3.65$$

$$\text{Higher fence} = 17.5 + 1.5 \times 4.5 = 24.25$$

$$[-3.65 \longleftrightarrow 24.25]$$

This means in the dataset 27 is an outlier.

$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27\}$

①  $\text{Min} = 1$

②  $Q_1 = 3$

③  $\text{Median} = 5$

④  $Q_3 = 7.5$

⑤  $\text{Max} = 9$

